

STATISTICAL INFERENCES IN CROSS-LAGGED PANEL STUDIES

TECHNICAL REPORT NO. 15

LAWRENCE S. MAYER

NOVEMBER 1985

U. S. ARMY RESEARCH OFFICE
CONTRACT DAAC29-85-K-0239

THEODORE W. ANDERSON, PROJECT DIRECTOR

DEPARTMENT OF STATISTICS
STANFORD UNIVERSITY
STANFORD, CALIFORNIA



APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED.

STATISTICAL INFERENCES IN CROSS-LAGGED PANEL STUDIES

Technical Report No. 15

Lawrence S. Mayer
Arizona State University and Stanford University

November 1985

U. S. Army Research Office
Contract DAAG29-85-K-0239

Theodore W. Anderson, Project Director

DEPARTMENT OF STATISTICS
STANFORD UNIVERSITY
STANFORD, CALIFORNIA

APPROVED FOR PUBLIC RELEASE: DISTRIBUTION UNLIMITED.

THE VIEW, OPINIONS, AND/OR FINDINGS CONTAINED IN THIS REPORT ARE THOSE OF THE AUTHOR(S) AND SHOULD NOT BE CONSTRUED AS AN OFFICIAL DEPARTMENT OF THE ARMY POSITION, POLICY, OR DECISION, UNLESS SO DESIGNATED BY OTHER DOCUMENTATION.

STATISTICAL INFERENCES IN CROSS-LAGGED PANEL STUDIES

1. Introduction

Panel studies are statistical studies in which two or more variables are observed for two or more subjects at two or more waves (points in time). In most panel studies the number of variables and the number of waves is small but the number of observations at any given wave is large. Cross-lagged panel studies are those studies for which the variables are continuous and the purpose of the study is to examine the cross-effects which are the impact of one set of variables on another over time. Such studies have been used in a variety of social and behavioral studies including studies which examine such controversial issues as the effect of IQ on achievement for school children and the effect of expenditures for police on a city's crime rate (e.g. Crano, Kenny, and Campbell (1972), Greenberg, Kessler, and Logan (1979), and Eaton (1978)). The purpose of this paper is to contribute to the statistical methods used in such studies.

Methods for analyzing cross-lagged panel studies have been developed over the past 20 years. Early methods were guided by the work of Donald Campbell (1963) and were motivated, to some degree, by the seminal work

*The author acknowledges the aid of his colleagues and students with particular gratitude to T.W. Anderson, D.R. Rogosa, S.S. Carroll and B.W. Brown. This research was conducted while the author was a Visiting Scholar, Department of Statistics, Stanford University.

of Paul Lazarsfeld (1948) on analyzing panel studies involving discrete data. The methods proposed by Campbell and developed by him and his colleagues (e.g. Cook and Campbell (1979) and Kenny (1979)) focused on the logic of using the correlational structure obtained from panel data to examine for the presence of cross-effects. For the most part sampling variation was ignored. Later authors argued that the presence of cross-effects was best examined by using a regression approach which involves formulating regression models of the responses, interpreting the cross-effects in terms of regression parameters, and using standard statistical methods to make inferences about the regression parameters and thus about the presence of cross-effects (e.g. Pelz and Andrews (1964), Duncan (1969), Heise (1970) and Rogosa (1980)). Kessler and Greenberg (1981) have provided an excellent review of the development of cross-lagged panel methods.

We contribute to the regression approach to the analysis of cross-lagged panel studies by examining issues that affect the optimality of the methods used to estimate the regression parameters and to test hypotheses about the presence of cross-effects. The first issue is the simultaneous nature of the regression models which arises from the fact that the regression approach can allow correlation among errors associated with different equations. The second is the assumption made on the observations of the initial wave, they can be fixed or observational in nature. Direct application of regression methods to panel studies often ignores the possibility of stochastic behavior for the initial observations. The third is the fact that the regression parameters may not be homogeneous in the sense that they change over waves.

Our approach is to define a statistical model for the panel study and then to consider the problems of estimation and testing for the parameters in the model. We borrow heavily from the theory of multivariate linear models (Anderson (1984), Rao (1973), and Duntzman (1984)) and from the theory of replicated vector-valued autoregressive processes (Anderson (1978)).

In the next section we introduce the model and set the notation. In the third section we present our results for the two-wave panel study and in the fourth section we extend these results to the multi-wave panel. We then apply our results to a panel study of the attitudes and perceptions of patients in a health maintenance organization and conclude with remarks on our current efforts.

2. Specification of the Model

Let $\mathbf{z}_{it} = (z_{it}^{(1)}, \dots, z_{it}^{(k)})'$ be k variables measured on the i th of n independent subjects at wave t ($t = 0, \dots, T$) of a panel study.

Suppose the variables divide into two sets as indicated by the partition $\mathbf{z}_{it} = (\mathbf{x}'_{it}; \mathbf{y}'_{it})'$ where \mathbf{x}_{it} and \mathbf{y}_{it} are of dimensions p and q respectively and $p + q = k$. The emphasis of the study is to estimate and test the effects of $\mathbf{x}_{i,t-1}$ and $\mathbf{y}_{i,t-1}$ on \mathbf{x}_{it} and \mathbf{y}_{it} .

The multivariate regression structure is

$$\mathbf{z}_{it} = \mathbf{B} \mathbf{z}_{i,t-1} + \boldsymbol{\varepsilon}_{it} \quad (i = 1, \dots, n; t = 1, \dots, T) \quad (2.1)$$

where the unobserved error vectors $\{\xi_{it}\}$ are independent and identically distributed random vectors with mean Q and covariance matrix \hat{A} , and B is a matrix of unknown regression coefficients.

The vector ξ_{it} and matrices $B = (b_{ij})$ and $A = (\sigma_{ij})$ are partitioned conformally with ζ_{it} to give

$$\xi_{it} = \begin{bmatrix} \xi_{it}^{(1)} \\ \xi_{it}^{(2)} \end{bmatrix} \quad B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} \hat{A}_{11} & \hat{A}_{12} \\ \hat{A}_{21} & \hat{A}_{22} \end{bmatrix}$$

In matrix notation the model is summarized as

$$\zeta_t = \zeta_{t-1} B' + E_t \quad (t = 1, \dots, T)$$

where $\zeta_t = (X_t; Y_t)$ has i th row ζ_{it} , $E_t = (E_t^{(1)}; E_t^{(2)})$ has i th row ξ_{it} , and the covariance matrix of E_t can be written $I_n \otimes \hat{A}$ where \otimes is the Kronecker product of linear algebra.

Specification of the model is completed by specifying the behavior of the observation matrix of the initial wave, ζ_0 . We either assume that ζ_0 is fixed or assume that the rows of ζ_0 are independent normal random vectors with common mean μ_0 and common covariance matrix θ_0 . If ζ_0 is fixed we let $Q_n = \frac{1}{n} \sum_{i=1}^n \zeta_{i0} \zeta_{i0}'$ and assume $Q = \lim Q_n$ exists and is of full rank.

If ζ_0 is random then the covariance matrix θ_1 of ζ_{i1} satisfies

$$\theta_1 = B \theta_0 B' + \hat{A} \quad (2.2)$$

To highlight the multivariate nature of the model we consider the simple case of $p = q = 1$ and $T = 1$. The model in (2.1) reduces to

$$x_{il} = b_{11}x_{i0} + b_{12}y_{i0} + \epsilon_{il}^{(1)} \quad (2.3)$$

$$y_{il} = b_{21}x_{i0} + b_{22}y_{i0} + \epsilon_{il}^{(2)}$$

The equations look like ordinary multiple regression equations with the added condition that the errors $\epsilon_{it}^{(1)}$ and $\epsilon_{it}^{(2)}$ are allowed to be correlated.

3. Two-wave Panel Studies

We consider the problems of parameter estimation and hypothesis testing in the two-wave panel study. Such consideration is pertinent because published applications are often two-wave studies (e.g. Crano, Kenny and Campbell (1979)) and methodological papers have been devoted to the statistical analysis of such studies (e.g. Duncan (1969)). Furthermore the results for two-wave studies are somewhat simpler than the results for multi-wave studies and thus "set the stage" for more difficult results.

We consider the problem of estimating the regression parameter matrix \hat{B} and covariance matrix $\hat{\Sigma}$ and then consider the problem of assessing the size of the effects and cross-effects. Technical aspects of the results are discussed briefly in an appendix.

The analysis of the two-wave model begins by defining two estimators

$$\hat{B} = (z_0' z_0)^{-1} z_0' z_1 = \left[\sum_{i=1}^n z_{i0} z_{i0}' \right]^{-1} \sum_{i=1}^n z_{i0} z_{i1}' \quad (3.1)$$

$$\hat{\Sigma} = \frac{1}{n} (z_1 - z_0 \hat{B}')' (z_1 - z_0 \hat{B}') = \frac{1}{n} \sum_{i=1}^n (z_{i1} - \hat{B} z_{i0}) (z_{i1} - \hat{B} z_{i0})'$$

The vector \hat{b}_j , the j th row of $\hat{\beta}$, contains the regression coefficients which can be obtained by regressing $z_{i1}^{(j)}$, the j th variable at wave 1, on $z_{i0}^{(1)}, \dots, z_{i0}^{(k)}$, all k variables at wave 0, using a standard multiple regression routine. The covariance estimator $\hat{\Sigma}$ cannot be obtained directly from the multiple regression output but can be computed using $\hat{\beta}$ and the numerical calculation indicated in (3.1).

3.1 Estimation of the Effects and Covariances

The use of the estimators in (3.1) can be justified, in the case of fixed initial observations, by applying results of multivariate linear models. We begin with

Theorem 1: For a two-wave panel study with fixed initial observations and normal errors the estimators $\hat{\beta}$ and $\hat{\Sigma}$ are maximum likelihood, \hat{b}_j , the j th row of $\hat{\beta}$, is a minimum variance unbiased estimator of b_j , the j th row of β , and the vectors $\{\sqrt{n}(\hat{b}_j - b_j)\}$ are normal random vectors with mean 0 and covariance matrices

$$\begin{aligned}\text{Var}(b_j) &= \sigma_{jj} Q_n^{-1} \\ \text{Cov}(b_j, b_k) &= \sigma_{jk} Q_n^{-1} \quad (j \neq k)\end{aligned}\tag{3.2}$$

This result follows directly from the theory of multivariate regression analysis (e.g. (Anderson (1984), Rao (1973))). See the appendix for additional comments.

If the assumption of normal errors is relaxed then Theorem 1 is replaced by

Theorem 2: For a two-wave panel model with fixed initial observations the estimator $\hat{\beta}_j$ of β_j ($j = 1, \dots, k$) is a minimum variance linear unbiased estimator and $\sqrt{n}(\hat{\beta}_j - \beta_j)$ is, for large n , approximately a normal random vector with the mean and variance given in theorem 1 with Q_n replaced by $Q = \lim_{n \rightarrow \infty} Q_n$.

This result, like Theorem 1, follows from applying the theory of multivariate linear models.

The covariance matrices of $\sqrt{n}(\hat{\beta}_j - \beta_j)$ in Theorems 1 and 2 must be calculated with the formulae indicated if a multiple regression routine is used to estimate the model but are provided as direct output if a multivariate regression routine, such as the routine in SAS, is used.

To this point we have assumed that the initial observation matrix Z_0 is fixed. Often panel studies are completely observational in nature and it may be more realistic to assume that Z_0 is random. For this case we have

Theorem 3: For a two-wave panel model with normal initial observations and normal errors, the estimators $\hat{\beta}$ and $\hat{\delta}$ given in (3.1) are maximum likelihood and the vectors $\{\sqrt{n}(\hat{\beta}_j - \beta_j)\}$ are for large n , approximately normal random vectors with mean Q and covariance matrices

$$\text{cov}[\sqrt{n}(\hat{b}_j - b_j)] = \sigma_{jj}\theta_0 \quad (3.3)$$

$$\text{cov}[\sqrt{n}(\hat{b}_j - b_j), \sqrt{n}(\hat{b}_l - b_l)] = \sigma_{jl}\theta_0 \quad (j \neq l)$$

Furthermore, the large-sample distribution of $\sqrt{n}(\hat{b}_j - b_j)$ remains the same if the condition of normal errors is relaxed.

The proof follows from the theory of replicated vector-valued auto-regressive processes (Anderson (1978)) and is discussed in the appendix.

Although the estimators given in (3.1) are maximum likelihood estimators for the model with normal initial observations and normal errors, the model does contain an additional vector μ_0 and an additional covariance matrix θ_0 , and both must be estimated. They are the mean vector and covariance matrix of the initial observation vector z_{i0} .

The maximum likelihood estimators are

$$\hat{\mu}_0 = \frac{1}{n} \sum_{i=1}^n z_{i0} / n \text{ and } \hat{\theta}_0 = \frac{1}{n} (z_{i0} - \hat{\mu}_0)(z_{i0} - \hat{\mu}_0)' / n \quad (3.4)$$

For the model with random initial observations computation of $\hat{\mu}_0$ can be done by a multiple regression routine as was true for the model with fixed initial observations, but computation of the maximum likelihood estimators of β , μ_0 , and θ_0 would require numerical calculation with the formulae in (3.1) and (3.4).

Similarly the covariance matrix of $\sqrt{n}(\hat{b}_j - b_j)$ can be computed by numerical calculation using the formula in (3.2).

If a multivariate regression routine is used then the estimators of β_j and the covariance matrix of $\sqrt{n}(\hat{\beta}_j - \beta_j)$ are given as output; the estimators of μ_0 and θ_0 must be evaluated by numerical calculation using the formula in (3.4).

Thus far we have assumed that the covariance matrix of the initial wave is different than the covariance matrix of the first wave. Alternatively, we could adopt an assumption common in econometrics, that the process has a long but unobserved history and thus has become stationary by wave 0. Then ζ_{10} and ζ_{11} would have identical covariance matrices, formally $\theta_0 = \theta_1$. This identity and the structure of the model imply that $\theta_1 = R\theta_0 R' + \Sigma$ and thus the covariance matrix θ_1 would satisfy

$$\theta_1 = R\theta_0 R' + \Sigma \quad (3.5)$$

Incorporating expression (3.5) as a constraint in the maximum likelihood optimization involves maximizing an equation which is highly non-linear in the elements of R . Optimization requires a maximum likelihood routine capable of handling nonlinear constraints. In our experience the values of the estimators obtained have been close to the values of the estimators given by (3.1). In theory, however, these estimators are not equivalent to the ones given in (3.1), even for large samples, nor is their behavior easily tractable. (see Anderson (1978) for similar remarks).

In addition to the computational complexity the stationary model is unattractive because panel data with only a few (4 or less) waves are not likely, in our experience, to be stationary. We prefer to treat θ_0 as an independent parameter to be estimated from the data.

3.2 Assessing the Overall Effects and the Cross-effects

Having considered the problem of estimating the regression parameters and covariance matrices of the two-wave model we turn to the problems of assessing the strength of the effects. We begin with consideration of the problem of testing the significance of the overall effect of z_{i0} , the variables at wave 0, on z_{il} the variables at wave 1. We then consider the problem of testing the significance of the cross-effects or effects of z_{i0} and z_{il} and of z_{i0} on z_{il} . We conclude with consideration of the problem of summarizing the strength of association due to the existence of cross-effects in the model.

For the model with normal errors the significance of the overall effect of z_{i0} on z_{il} is tested by any of several standard procedures used to test the hypothesis $H_0^{(1)}$: $\beta = 0$ in multivariate linear models. These tests are covered in texts on multivariate statistics (Anderson (1984), Rao (1973)) and are evaluated by most standard multivariate regression programs such as the routine in SAS.

Of the various tests (likelihood ratio, Lawley-Hotelling, Pillai, etc.) we tend to prefer the likelihood ratio test, in part, because it is the most theoretically compatible with maximum likelihood estimation and the latter is used to estimate the models.

For many panel studies the overall effects of the variables at wave 0 on the variables at wave 1 are highly significant owing, in part, to the high degree of serial correlation found in each individual variable. Formally, the diagonal elements of β , the effects of $z_{i0}^{(i)}$ on $z_{il}^{(i)}$, tend to be large (see the example in section 5).

The more tentative and more intriguing issue is the significance of the effects of each subset of variables at wave 0 on the other subset of 1, the effects of x_{i0} on x_{i1} and of x_{i0} on x_{il} . These cross-effects are the off-diagonal blocks, β_{12} and β_{21} , of β and the hypothesis to be tested, the hypothesis of no cross-effects is $H_0^{(2)}$: $\beta_{12} = \beta'_{21} = 0$.

The hypothesis of no cross-effects cannot be tested with a standard multivariate regression routine because the null hypothesis $H_0^{(2)}$ cannot be expressed in terms of linear contrasts in β (see the appendix for additional comments).

One approach to testing $H_0^{(2)}$ is to divide the hypothesis into $\beta_{12} = 0$ and $\beta_{21} = 0$ and then apply standard statistical theory to test each of the component hypotheses. The components can be represented as standard contrasts in β and thus a multivariate regression routine yields the test. Then some method is used to combine results of the two tests into a test of $H_0^{(2)}$. One method (Kessler and Greenberg (1981)) combines the two tests by ignoring the correlation between them, and thus, implicitly assumes the tests are independent. This method seems reasonable, but may give an advertised probability of type 1 error quite different from the true probability of a false rejection, the difference arising from the correlation between the errors of the equations.

We choose not to divide $H_0^{(2)}$ but, instead, to develop the likelihood ratio test of $H_0^{(2)}$, a test which takes into account the correlation structure of the errors. We are not able to give a closed form expression for the test statistic but instead, evaluate the test statistic directly from two applications of a standard maximum

likelihood estimation routine. The easiest explanation of the likelihood ratio test involves presenting the model in the widely used notation associated with the LISREL routine for maximum likelihood estimation (Jöreskog (1979)). We will keep notations distinct by using a superscript dot to indicate vectors and matrices in the LISREL notation.

Begin with the LISREL model

$$\dot{y} = \dot{\Lambda}_y \eta + \dot{\xi} \quad (3.5)$$

$$\dot{x} = \dot{\Lambda}_x \xi + \dot{\delta} \quad (3.6)$$

$$\dot{\eta} = \dot{B} \eta + \dot{\Gamma} \xi + \dot{\zeta} \quad (3.7)$$

where \dot{y} and \dot{x} are observed random vectors; η , ξ , and ζ are unobserved random variables; $\dot{\Lambda}_y$, $\dot{\Lambda}_x$, \dot{B} , and $\dot{\Gamma}$ are unknown parameter matrices; and

$$\dot{\theta}_\epsilon = \text{cov}(\dot{\xi}) \quad \dot{\theta}_\delta = \text{cov}(\dot{\delta})$$

$$\dot{\phi} = \text{cov}(\dot{\xi}) \quad \dot{\psi} = \text{cov}(\dot{\zeta})$$

To represent the two-wave panel model in this notation let $\dot{y} = (z'_{i1}, \dots, z'_{nl})'$ and $\dot{x} = (z'_{i0}, \dots, z'_{n0})'$ be nk dimensional vectors and let $\dot{\Lambda}_x = \dot{\Lambda}_y = 0$ and $\dot{\theta}_\epsilon = \dot{\theta}_\delta = 0$. Equations (3.5) and (3.6) become trivial identities and (3.7) becomes

$$\dot{y} = \dot{B} \dot{y} + \dot{\Gamma} \dot{x} + \dot{\zeta} \quad (3.8)$$

Assuming $\dot{B} = 0$, $\dot{\Gamma} = B$, $\dot{\phi} = \theta_0$, and $\dot{\psi} = \psi$ completes the representation of the two-wave model into LISREL notation.

Let the model specified above be fit to the data and let χ^2 be the chi-square statistic produced by LISREL which indicates the badness-of-fit of the model.

To obtain the likelihood ratio test of $H_0^{(2)}$ a second model is fit to the data, a model which is the same as the model above except that the regression parameter matrix B is constrained to be block diagonal as indicated by the null hypothesis $H_0^{(2)}: R_{12} = R'_{21} = Q$. Let χ_0^2 be the chi-square badness of fit statistic for this constrained model.

The likelihood ratio test statistic is the simple difference $t_1 = \chi_0^2 - \chi^2$ and the test of the hypothesis of no cross-effects proceeds by treating t_1 as a chi-square statistic with $2pq$ degrees of freedom.

Note that the test presented above requires two maximizations of the likelihood function, or equivalently, two runs of LISREL. Two alternatives exist, neither of which requires these two maximizations and both of which can be performed with a multivariate regression routine.

The first alternative is to approximate the test statistic t_1 by application of Zellner's theory of seemingly unrelated regressions (Zellner (1962)). This approximation is outlined in the appendix.

The second alternative is to use a test based on the distribution of the estimators of the cross-effects ignoring the null hypothesis. Such tests are called Wald-type tests (eg. Judge, et. al. (1980)). The test

statistic, t_2 , is obtained from the asymptotic distribution of $\sqrt{n}(\hat{b}_{v_j} - b_{v_j}')$. In particular we let $\hat{b}^{(vec)} = (\hat{b}_1', \dots, \hat{b}_k')'$ be the vector formed from the rows of \hat{B} and $b^{(vec)}$ the vector formed from the rows of B . Let U be the $k^2 \times k^2$ covariance matrix of $\sqrt{n}(\hat{b}^{(vec)} - b^{(vec)})$ pieced together from the covariance matrices given for $\sqrt{n}(\hat{b}_{v_j} - b_{v_j}')$ in Theorems 1, 2 or 3. Let $\hat{b}_{v12}^{(vec)}$ be the subvector of $\hat{b}^{(vec)}$ containing the elements of \hat{B}_{v12} and \hat{B}_{v21} and let $b_{v12}^{(vec)}$ be the corresponding subvector of $b^{(vec)}$; let U_{12} be the $2pq \times 2pq$ submatrix of U containing the covariance matrix of $\sqrt{n}(\hat{b}_{v12}^{(vec)} - b_{v12}^{(vec)})$. Under the hypothesis of no cross-effects $b_{v12}^{(vec)} = 0$ and the test statistic is

$$t_2 = \hat{b}_{v12}^{(vec)'} U_{12}^{-1} \hat{b}_{v12}^{(vec)}$$

which is for large samples, a chi-square statistic with $2pq$ degrees of freedom.

One advantage of the test based on t_2 over the test based on t_1 is that if the assumption of normal errors is relaxed the test based on t_2 remains accurate for large samples.

If the assumption of normal errors holds then the tests based on the likelihood ratio statistic t_1 and the Wald-type statistic t_2 are almost identical for large samples. For samples of moderate size the powers of the tests can be compared by selecting an alternative hypothesis

$H_A: B_{v12} = B_A$ and $B_{v21} = B_B$, generating multiple data sets with B_A and B_B as the true values of B_{v12} and B_{v21} , and then calculating the proportion of data sets for which each test rejects. Preliminary calculations

suggest that the powers of the tests are very close if n is 50 or larger. Note however that such simulation gives indication of the performance of the tests when the errors are truly normal. We conjecture the Wald-type test may perform slightly better than the likelihood test if the errors are far from normal. Obviously, a more complete simulation study is needed (see Evans and Savin (1982) for more on the relationship between these tests and other tests when used with an econometric model containing lagged dependent variables; Also, see Rothenberg (1982) for the argument that for the multiple regression model the tests have similar power properties for samples of moderate sizes. These results do not apply directly to the cross-lagged panel model - because of the presence of lagged predictors - but suggest further study is needed.

We close our analysis of the two-wave model by considering the problem of summarizing the degree of association in the model due to the inclusion of the cross-effects.

One of the most widely used statistics in interpreting the output of a multiple regression analysis is the square of the partial correlation coefficient which indicates the percentage of variation in the dependent variable explained by some variables controlling for others. This statistic is particularly attractive in that it is a proportional reduction in error statistic where the sum of squares residual is used as the measure of error. Thus it has a simple interpretation which, at times, is valid across models and variables.

Following Sobel and Boernstedt (1985) a proportional reduction in error measure somewhat similar to the squared partial correlation coefficient can be used to summarize the reduction in the chi-square badness-of-fit statistic when the cross-effects are included in this model. To be specific we suggest the measure

$$PRE = (\chi^2_0 - \chi^2) / \chi^2_0$$

gives some indication of the importance of the cross-effects in the panel model controlling for the other effects. Note that Sobel and Boernstedt extend this type of measure to one that compares the models with and without cross-effects to "baseline models." These extended measures could be used with the two-wave panel model. The reader should consult their work for details.

4. Multi-wave Panel Studies

The issues that arise in estimating the parameters and testing the hypothesis of no cross-effects for a multi-wave study are somewhat different than those found in the analysis of a two-wave study. First, the issue of whether to treat the initial wave as fixed or random is less critical since the regressor variables include lagged endogenous variables under either assumption. Secondly, the assumption that B is homogenous, or constant, over waves is often questionable and is open to examination.

For a study of $T + 1$ waves ($T = 2, 3, \dots$) the regression structure in (2.1) extends to

$$\bar{z}_{it} = B \bar{z}_{i,t-1} + \xi_{it} \quad i = 1, \dots, n \quad t = 1, \dots, T \quad (4.1)$$

If we $\mathbf{z}_t = (z_{1t}, \dots, z_{nt})'$ and $\mathbf{e}_t = (e_{1t}, \dots, e_{nt})'$ then the regression structure in (4.1) can be summarized as

$$\mathbf{z}_t = \mathbf{z}_{t-1} \hat{\mathbf{B}}' + \mathbf{e}_t$$

We assume $(\sum_{t=0}^{T-1} \mathbf{z}_t' \mathbf{z}_t)$ is of full rank and define two estimators

$$\hat{\mathbf{B}} = (\sum_{t=0}^{T-1} \mathbf{z}_t' \mathbf{z}_t)^{-1} (\sum_{t=1}^T \mathbf{z}_{t-1} \mathbf{z}_t')$$

and

(4.2)

$$\hat{\mathbf{x}} = (Tn)^{-1} \sum_{t=1}^T (\mathbf{z}_t - \mathbf{z}_{t-1} \hat{\mathbf{B}}') (\mathbf{z}_t - \mathbf{z}_{t-1} \hat{\mathbf{B}}')'$$

These estimators are pooled least-squares estimators and are natural extensions of the estimators given in (3.1).

The vector $\hat{\mathbf{b}}_j$, the j th row of $\hat{\mathbf{B}}$ can be obtained by regressing all Tn observations of the form $z_{it}^{(j)}$ on the appropriate values of $z_{1,t-1}^{(1)}, \dots, z_{i,t-1}^{(k)}$ with a multiple regression routine. Then $\hat{\mathbf{x}}$ can be computed from (4.2). Alternatively $\hat{\mathbf{B}}$ and $\hat{\mathbf{x}}$ can be obtained directly as output of a multivariate regression routine, or of a maximum likelihood estimation routine such as LISREL. As a partial extension of Theorem 1 we have

Theorem 4: For the multi-wave panel study with normal errors the estimators $\hat{\mathbf{B}}$ and $\hat{\mathbf{x}}$ given in (4.2) are maximum likelihood whether the initial observations are fixed or normal.

The exact distribution of $\hat{\mathbf{B}}$ does not follow from the normality of the errors as it did in the two-wave model, and thus no extension of the

distribution result in Theorem 1 is possible. Theorems 2 and 3 do extend, in part, to

Theorem 5: For the multi-wave panel study with fixed initial observations $\hat{\beta}$ and $\hat{\Sigma}$ are consistent estimators (in n) and the vectors $\{\sqrt{n}(\hat{b}_j - b_j)\}$ are, for large samples, normal random vectors with mean θ and covariance matrices

$$\text{cov}[\sqrt{n}(\hat{b}_j - b_j)] = \sigma_{jj} Q^{-1}$$

$$\text{cov}[\sqrt{n}(\hat{b}_j - b_j), \sqrt{n}(\hat{b}_\ell - b_\ell)] = \sigma_{j\ell} Q^{-1}$$

If the initial observations are normal then $\{\sqrt{n}(\hat{b}_j - b_j)\}$ have the same large-sample distribution with Q^{-1} replaced by θ_0 .

We turn to the problem of testing the significance of the effects and cross-effects in the multi-wave panel.

We again adopt the notation of the LISREL routine. We use the identifications of section 3 with one extension which is as follows for the 3-wave model:

$$\dot{x} = (z'_{11}, \dots, z'_{nl}; z'_{12}, \dots, z'_{n2})' \quad \dot{x} = (z'_{10}, \dots, z'_{n0})'$$

and let

$$\begin{aligned} \dot{\beta} &= \begin{bmatrix} \theta & \theta \\ \theta & \theta \end{bmatrix} & \dot{\Sigma} &= \begin{bmatrix} \hat{\beta} \\ \theta \end{bmatrix} \\ \dot{\psi} &= \theta_0 & \dot{\psi} &= \begin{bmatrix} \hat{\beta} & \theta \\ \theta & \hat{\beta} \end{bmatrix} \end{aligned}$$

The representation for a general $T + 1$ wave panel is a direct extension of the above.

The likelihood ratio test of the significance of the effects, or equivalently, the test of hypothesis $\beta = Q$ is not produced by either a multiple regression routine or a multivariate regression routine. It is obtained from LISREL by fitting the model above to the data and then fitting the model again but with the constraint $\beta = Q$. The test statistic, t , is the difference between the chi-square badness-of-fit statistics for the two models. For large samples it is approximately a chi-square random variable with $k^2 = (p + q)^2$ degrees of freedom.

Turning to the problem of testing the significance of the cross-effects in the multi-wave panel models, the likelihood ratio test is generated by a direct extension of the one used for the two-wave model. Let $\chi^2(T + 1)$ be the chi-square badness-of-fit measure for the $T + 1$ wave model and let $\chi_0^2(T + 1)$ be the same measure under the constraint that $\beta_{12} = \beta'_{21} = Q$, or equivalently, under the hypothesis that β is block diagonal. The test statistic is the difference

$$t_3 = \chi_0^2(T + 1) - \chi^2(T + 1)$$

which for large samples is approximately a chi-square random variable with $2pq$ degrees of freedom.

If a multiple regression routine or a multivariate regression routine is used then the likelihood ratio test can be approximated by using a method similar to the method used to approximate the test statistic t_1 in section 3.

A Wald-type test of the significance of the cross-effects based on the large sample distribution of $\hat{\beta}$ (for the multi-wave model) can also be used. The test statistic t_4 is the natural generalization of the test based on t_2 for the two-wave model and is summarized in the appendix.

The test based on t_4 , like the test based on t_2 for the two-wave panel, can be used with either the assumption of normal initial observations or the assumption of normal errors relaxed.

Note that the hypothesis of no cross-effects, as we have defined it, is not the hypothesis that the sets of variables X_t and Y_t are independent. The hypothesis of independence is stronger than the hypothesis $H_0^{(2)}$.

Assuming the errors are normal this stronger hypothesis is the intersection of the hypothesis of no cross-effects and the hypothesis

$$H_0^*: \begin{matrix} \ddot{t}_{12} = \ddot{t}'_{21} = 0 \end{matrix} \text{ where } \ddot{t}_{12} \text{ and } \ddot{t}'_{21} \text{ are the off-diagonal blocks of } \ddot{\beta}.$$

Anderson (1978) develops a test for the hypothesis H_0^* and then a conditional test for the hypothesis of no cross-effects given H_0^* is true. We do not use his test but note that it can be applied directly to the multi-wave panel model with normal initial observations and normal errors.

Finally, for a multi-wave panel study the assumption of homogeneity of the parameter matrix β across can be tested. Relaxing this assumption the model in (4.1) becomes

$$Z_t = Z_{t-1} \beta_t' + E_t \quad t = 1, \dots, T$$

where the terms are as defined in Section 2 except that the parameter matrix β_t depends on t .

The maximum likelihood estimator of β_t is

$$\hat{\beta}_t = (\hat{Z}'_{t-1} \hat{Z}_{t-1})^{-1} \hat{Z}'_{t-1} \hat{Z}_t$$

In the context of the multivariate autoregressive process, Anderson (1978) develops a test of the hypothesis of homogeneity $H_0: \beta_1 = \dots = \beta_T$ which uses the test statistic.

$$t_5 = n \operatorname{tr} \sum_{t=1}^T (\hat{\beta}_t - \hat{\beta}) [\frac{1}{n} \hat{Z}'_{t-1} \hat{Z}_{t-1}] (\hat{\beta}_t - \hat{\beta})' \hat{\beta}^{-1}$$

where $\hat{\beta}$ and $\hat{\beta}$ are defined in (4.2). If the initial observations are normal and the errors are normal then t_5 is, for large samples, almost a chi-square random variable with $(T-1)pq$ degrees of freedom. For multi-wave panel studies t_5 can be used to test the assumption of homogeneity.

The Sobel and Bohrnstedt proportional reduction in error measure extends immediately from the two-wave to a multi-wave model.

5. Application of Methods to a Panel Study

The upper management of a consortium of health maintenance organizations (HMO's) wants to know if there is significant correlation between patients attitudes toward health maintenance organizations and their perceptions of the quality of care they are receiving from the HMO in which they enrolled. If so, they would like to know if over time there

appears to be a "causal priority" between such attitudes and perceptions. Does the patient's attitudes toward HMO's precede or "drive" his or her perceptions of the quality of the care being received? If yes, then management might want to invest resources in a campaign to improve the attitudes of the general public toward HMO's. On the other hand suppose the effect over time is reversed with the perceptions of the quality of care preceding or "driving" the attitude toward HMO's. Then management may want to invest those same resources in a more focused campaign to improve the patient's perceptions of the care received.

To obtain a preliminary insight into the issue, management conducted a survey of randomly selected patients enrolled in their member HMO's. For a variety of reasons, including minimizing cost and disruption, a panel design was used. Patients were interviewed upon completion of each self-initiated visit to the HMO, the visits being considered waves. To demonstrate our methods we analyze a subsample of 50 patients over three waves. This subsample was chosen randomly but for the sake of simplicity patients with incomplete data or unusual health-care status (eg. the terminally ill) were not considered for this subsample.

Two indicators of attitudes toward HMO's are used in our analysis. The first, X_1 , indicates the patients attitude toward the specific HMO in which or she is enrolled. The second X_2 , indicates the patients attitude toward the concept of "socialized medicine" meaning the government providing the general public a low-cost alternative to fee-for-service health-care. These variables were thought to capture two closely related dimensions in the overall attitude toward HMO's. Both

variables are rescaled, in part to mask propriety data, to have mean 10 and a standard deviation of about 3 with a higher value indicating a more positive attitude.

Two perceptions of quality of health-care received are included. The first, Y_1 , indicates the perceptions of the quality of care received in the visit just concluded and the second, Y_2 , indicates the perception of the quality of care received since initial enrollment. The two variables, while related, were thought to capture different issues of quality of care. Again, the variables were scaled over a larger sample to have mean of 10 and standard deviation of about 3 with a higher value indicating a more positive perception.

Table 1 displays some correlational statistics for the raw data. The correlation matrix shows that there is some (marginal) relationship between the two indicators of the attitude toward HMO's ($r = .38$) but little relationship between the two pairs of measurements. The multiple correlations indicate that the attitudes and perceptions at wave t are well predicted by the attitudes and perceptions at wave t - 1. The comparison between the multiple correlation and serial correlations indicates that the majority of this predictability can be attributed to the relationship between a measurement at wave t - 1 and the same measurement at wave t. For each variable the serial correlation is within .02 of the multiple correlation. Along with these summary statistics the usual data descriptive methods (histograms, box plot, stem and leaf plots, and others) were used to examine the shapes of the distributions of the measurements. All were quite symmetric and fairly normal with no significant outliers.

Table 2 summarizes the output of a multiple regression routine applied to each variable individually. The estimates of the regression parameters are the maximum likelihood estimates of Theorem 4.

The cross-effects seem small when compared to their estimated (marginal) standard error but Table 2 does not provide a correct test of the significance of the effects or cross-effects. For these we turn to Table 3 which again displays the maximum likelihood estimates of the regression parameters. These estimates can be obtained from the regression output in Table 1, from a multivariate regression routine, or from a maximum likelihood routine such as LISREL.

The ratio of each estimate to its (correctly) estimated standard errors are given in Table 3. The estimates of the standard errors used as the denominators are obtained from a multivariate regression routine, from LISREL; or from the multiple regression output. These estimates are the estimates of Theorem 5 and are accurate for large samples. The ratios are approximately normal and thus give a rough indication of the size of the coefficient. These ratios indicate that attitudes at wave $t - 1$ are good predictors of attitudes at wave t and that perceptions at wave $t - 1$ are good predictors of perceptions at wave t .

The ratio of the estimates of the cross-effect parameters to their estimated standard errors are all quite small. These ratios indicate that neither attitudes nor perceptions are good predictors of the other over waves. We note their sizes but do not use these ratios to test the hypothesis of no cross effects since the ratios are not independent.

Table 3 also gives the likelihood ratio statistic for testing the hypothesis of no effects ($t = 688.69$) which can be compared to the critical value of a chi-square distribution with 16 degrees of freedom and is quite significant. This result is as expected given the high degree of serial correlation indicated in Tables 1 and 2.

The likelihood ratio test statistic for the hypothesis of no cross-effects is also given ($t_3 = 4.81$) as is the Wald-type test statistic ($t_4 = 5.41$). The former was obtained from two runs of the LISREL routine and the latter was obtained by brute calculation using the output of a multivariate regression routine. We found the PROC routines of SAS particularly handy for both the multivariate regression and for the subsequent calculations.

Both of these test statistics are asymptotically chi-square with 8 degrees of freedom and neither is significant at the .05 level. Neither test gives much evidence for the existence of cross-effects.

The final hypothesis of interest is the hypothesis of homogeneity, $B_{\cdot 1} = B_{\cdot 2}$. The test statistic for homogeneity ($t_5 = 173.11$) is significant at the .05 level. This test indicates that the matrix of regression parameters appears to vary over waves.

From the LISREL output the chi-square test for the goodness of fit for the model with cross-effects is $\chi^2 = 78.18$ and the value for the model without cross-effects is $\chi^2_0 = 82.96$. Thus the Sobel-Bohrnstedt type

proportional reduction in error measure has value (PRE) $\frac{82.96 - 78.18}{82.96} = .06$ which indicates that adding the cross-effects reduces the badness-of-fit of the model by about 6 percent.

The analysis suggests that the attitudes toward HMO's and the perceptions of the quality of health-care received are not highly interactive over waves. The conclusion is tentative, in part because the evidence against the homogeneity of the regression parameter makes interpretation difficult.

6. Future Directions

We have presented statistical results for estimating the parameters and testing the hypotheses of no effects and no cross-effects in a cross-lagged panel study. Our results are an extension of the work of several social methodologists on the regression approach to modeling panel data. They are also an application of results on multivariate linear inference applied to the panel models widely used by social scientists. We extend this earlier work by considering the simultaneous nature of the regression models formulated and by considering the nature of the observations made at the initial wave. We provide multivariate estimators of the regression parameters and tests of the hypotheses of no effects and cross-effects. We also consider the problems of summarizing the contributions of the cross-effects to the degree of fit of the model and the problem of testing the homogeneity assumption on the regression coefficients.

The continuous variable panel study has attracted attention in many disciplines including econometrics. There focus has been on the problem of correctly formulating the error structure in the regression approach and on the problem of estimating the model under different formulations (e.g. Balestra and Nerlove (1966), Maddala (1971) and Wallace and Hussain (1969)). Judge, et. al. (1980) give an excellent review of this research.

Anderson and Tsai (1981; 1982) have produced an excellent piece on some of the statistical issues that arise from a particular econometric specification of the error structure in a univariate panel model. They, and most econometricians, study a regression model for a single response variable in which the error for respondent i at time t can be decomposed into the sum of two (or more) terms, the first of which depends on i and t and the second of which depends only on i . The second term allows the model to capture the fact that the respondent may tend to be above the predicted value from the regression model at every wave. Various sets of assumptions can be made on the behavior of these two terms over.

For several panel studies we have adopted this econometric error formulation and it seems to be quite realistic. We are studying the problem of testing for the presence of cross-effects with this type of error structure (Mayer (1985b)).

We are also studying the problem of detection of autoregressive errors in the multivariate panel model as specified in this paper and the problem of testing for the presence of cross-effects given that the errors are autoregressive (Mayer (1985a)).

Finally we are looking at the problem of estimating the parameter matrix and covariance matrix and testing for the presence of cross-effects in a panel study where the homogeneity assumption on the matrix of regression parameters does not hold.

Table 1

Correlation Statistics for the Raw Data: HMO Data

Correlation Matrix for Raw Data

| | | | | |
|----------------|------|------|------|------|
| X ₁ | 1.00 | .38 | .20 | .20 |
| X ₂ | | 1.00 | .12 | .05 |
| Y ₁ | | | 1.00 | .38 |
| Y ₂ | | | | 1.00 |

Multiple Correlations and Serial Correlations

| Multiple Correlations | | Serial Correlations | |
|-----------------------|-----|---------------------|-----|
| X ₁ : | .96 | X ₁ : | .95 |
| X ₂ : | .91 | X ₂ : | .87 |
| Y ₁ : | .93 | Y ₁ : | .93 |
| Y ₂ : | .75 | Y ₂ : | .70 |

Table 2
Multiple Regression Summary: HMO Data

Dependent Variable: X_1 (attitude toward HMO)

| Source | df | Sum of Squares | Mean Square | F | Prob Value |
|--------|-----|----------------|-------------|-------|------------|
| Model | 4 | 1041.27 | 260.32 | 321.5 | < .0001 |
| Error | 96 | 77.73 | .81 | | |
| Total | 100 | 1119.00 | | | |

| Predictor Variable | df | Parameter Estimate | Standard Error | t | Prob Value |
|--------------------|----|--------------------|----------------|-------|------------|
| X_1 | 1 | .79 | .03 | 31.33 | < .0001 |
| X_2 | 1 | -.04 | .04 | -1.06 | .29 |
| Y_1 | 1 | -.03 | .03 | -1.05 | .30 |
| Y_2 | 1 | -.03 | .04 | -.56 | .57 |

Dependent Variable: X_2 (attitude toward socialized medicine)

| Source | df | Sum of Squares | Mean Square | F | Prob Value |
|--------|-----|----------------|-------------|--------|------------|
| Model | 4 | 513.73 | 128.43 | 116.12 | < .0001 |
| Error | 96 | 106.18 | 1.11 | | |
| Total | 100 | 619.90 | | | |

| Predictor Variable | df | Parameter Estimate | Standard Error | t | Prob Value |
|--------------------|----|--------------------|----------------|-------|------------|
| X_1 | 1 | -.06 | .03 | -1.97 | .05 |
| X_2 | 1 | .87 | .04 | 19.99 | .88 |
| Y_1 | 1 | .00 | .04 | .12 | < .0001 |
| Y_2 | 1 | -.02 | .05 | -.34 | .50 |

Table 2
(Continued)

Dependent Variable: Y_1 (perception of quality of current care)

| Source | df | Sum of Squares | Mean Square | F | Prob Value |
|--------|-----|----------------|-------------|--------|------------|
| Model | 4 | 443.85 | 110.96 | 137.01 | < .0001 |
| Error | 96 | 77.75 | .81 | | |
| Total | 100 | 521.59 | | | |

| Predictor Variable | df | Parameter Estimate | Standard Error | t | Prob Value |
|--------------------|----|--------------------|----------------|-------|------------|
| X_1 | 1 | .03 | .03 | 1.19 | .24 |
| X_2 | 1 | .01 | .04 | .15 | .88 |
| Y_1 | 1 | .65 | .03 | 20.59 | < .0001 |
| Y_2 | 1 | .03 | .04 | .67 | .50 |

Dependent Variable: Y_2 (perception of overall care)

| Source | df | Sum of Squares | Mean Square | F | Prob Value |
|--------|-----|----------------|-------------|-------|------------|
| Model | 4 | 274.41 | 68.60 | 30.70 | < .0001 |
| Error | 96 | 214.54 | 2.23 | | |
| Total | 100 | 488.96 | | | |

| Predictor Variable | df | Parameter Estimate | Standard Error | t | Prob Value |
|--------------------|----|--------------------|----------------|------|------------|
| X_1 | 1 | .01 | .04 | .19 | .85 |
| X_2 | 1 | .01 | .06 | .11 | .91 |
| Y_1 | 1 | .01 | .05 | .13 | .89 |
| Y_2 | 1 | .72 | .07 | 9.71 | < .0001 |

Table 3
Inferences for Panel Model: HMO Data

Estimates of Regression Parameters

| | | | | |
|----|------|------|-----|------|
| B: | .79 | -.04 | .03 | -.03 |
| v | -.06 | .87 | .00 | .02 |
| | .03 | .00 | .65 | .03 |
| | .01 | .01 | .01 | .72 |

Ratio of Estimates to Estimated Standard Errors

| | | | | |
|----|-------|-------|-------|------|
| B: | 31.32 | -1.06 | -1.05 | -.56 |
| v | -1.97 | 20.00 | .12 | -.34 |
| | 1.20 | .15 | 20.59 | .67 |
| | .19 | .11 | .13 | 9.71 |

Estimate of Covariance Matrices

| | | | | | | | | | |
|--------------------|-----|------|------|-----|--------------------|-------|-------|-------|------|
| $\hat{\alpha}_1$: | .91 | .01 | -.14 | .20 | $\hat{\theta}_0$: | 21.47 | 6.73 | 4.13 | 4.20 |
| | | 1.37 | | .03 | | | 8.393 | 1.81 | .82 |
| | | | | .28 | | | | 13.72 | 3.80 |
| | | | | .53 | .06 | | | | 5.95 |
| | | | | | 2.59 | | | | |

Likelihood Ratio Test of Presence of Effects

$t = 688.69$ approximately chi-square df = 16

Likelihood Ratio Test of Presence of Cross-Effects

$t_3 = 4.81$ approximately chi-square df = 8

Wald Type Test of Presence of Cross-Effects

$t_4 = 5.41$ approximately chi-square df = 8

Anderson Test of Homogeneity ($H_0: B_1 = B_2$)

$t_5 = 173.11$ approximately chi-square df = 4

Sobel-Bohrnstedt Proportional Reduction in Error Measure

PRE = .06

References

- Anderson, T.W. (1984) *An Introduction to Multivariate Statistical Analysis*, 2nd Edition, New York: Wiley.
- Anderson, T.W. (1978) "Repeated Measures on Autoregressive Processes," *Journal of the American Statistical Association*, 73, 371-378.
- Anderson, T.W. and Tsiao, C. (1981) "Estimation of Dynamic Models with Error Components," *Journal of the American Statistical Association*, 76, 598-606.
- Anderson, T.W. and Tsiao, C. (1982) "Formulation and Estimation of Dynamic Models Using Panel Data," *Journal of Econometrics*, 18, 47-82.
- Balestra, P. and Nerlove, M. (1966) "Pooling Cross-Section and Time Series Data in the Estimation of a Dynamic Model: The Demand of Natural Gas," *Econometrica*, 34, 585-612.
- Campbell, D.T. (1963) "From Description to Experimentation: Interpreting Trends as Quasi-experiments," in *Problems in the Measurement of Change* (C.W. Harris, ed), Madison: University of Wisconsin Press, 212-254.
- Campbell, D.T. and Stanley, J.C. (1963) *Experimental and Quasi-Experimental Designs for Research*, Chicago: Rand-McNally.
- Cook, T.D. and Campbell, D.T. (1979) *Quasi-Experimentation: Design and Analysis*, Chicago: Rand-McNally.

Cox, D.R. and Hinkley, D.V. (1974) *Theoretical Statistics*, London: Chapman and Hall.

Crano, W.D., Kenny, D.A. and Campbell, D.T. (1972) "Does Intelligence Cause Achievement: A Cross-Lagged Panel Analysis", *Journal of Educational Psychology*, 63, 258-275.

Duncan, O.D. (1969) "Some Linear Models for Two-Wave, Two-Variable Panel Analysis", *Psychological Bulletin*, 72, 177-182.

Dunteman, G. (1984) *Introduction to Multivariate Analysis*, Beverly Hills - Sage.

Eaton, W. (1978) "Life Events, Social Supports, and Psychiatric Symptoms", *Journal of Health and Social Behavior*, 19, 230-234.

Evans, G. and Savin, N. (1982) Conflict Among Testing Procedures in a Linear Regression Model with Lagged Dependent Variables, *Advances in Econometrics*, (W. Hildenbrand, ed.) Cambridge: Cambridge University Press.

Greenberg, D.G., Kessler, R.C. and Logan, C.H. (1979) "A Panel Model of Crime Rates and Arrest Rates", *American Sociological Review*, 44, 843-850.

Heise, D. (1970) "Causal Inference from Panel Data, Sociological Methodology", 1970 (E. Borgatta and G. Bohrnstedt, eds.) San Francisco: Jossey-Bass, 3-27.

- Joreskog, K. (1979) Statistical Models and Methods for the Analysis of Longitudinal Data, Chapter 5 in Advances in Factor Analysis and Structural Equation Models (K. Joreskog and D. Sorbom, eds.)
- Judge, G., Griffiths, W., Hill, A., Lee, T. (1980) The Theory and Practice of Econometrics, New York: Wiley.
- Kenny, D. (1979) Correlation and Causation, New York: Wiley.
- Kessler, R.C. and Greenberg, D.F. (1981) Linear Panel Analysis: Models of Quantitative Change, New York: Academic Press.
- Lazarsfeld, P.F. (1948) "The Use of Panels in Social Research", in Continuities in the Language of Social Research (P.F. Lazarsfeld, et.al. eds.) New York: Free Press, 330-337.
- Maddala, G.S. (1971) "The Use of Variance Components Models in Pooling Cross-Section and Time Series Data", Econometrica, 39, 341-358.
- Malinvaud, E. (1980) Statistical Methods of Econometrics, Amsterdam: North Holland.
- Mayer, L.S. (1985a) On Cross-Lagged Panel Models with Serially Correlated Errors, Technical Report, Department of Statistics, Stanford University.
- Mayer, L.S. (1985b) Multivariate Panel Models with Individual Effects in the Error Structure, Working Paper DIS 84/85-13, College of Business, Arizona State University.

Pelz, D.C. and Andrews, F.M. (1964) "Detecting Causal Priorities in Panel Study Data", American Sociological Review, 29, 836-848.

Rao, C.R. (1973) Linear Statistical Inference and Its Applications, Second Edition, New York: Wiley.

Rogosa, D.R. (1980) "A Critique for Cross-lagged Correlation", Psychological Bulletin, 88, 245-258.

Rothenberg, T. (1982) Comparing Alternative Asymptotically Equivalent Tests, in Advances in Econometrics (W. Hildenbrand, ed.) Cambridge: Cambridge University Press, p. 255-262.

Schmidt, P. (1976) Econometrics, New York: Dekker.

Sobel, M. and Bohrnstedt, G. (1985) Use of Null Models in Evaluating the Fit of Covariance Structure Models, in Sociological Methodology, 1985 (N. Tuma, ed.), San Francisco, Jossey-Bass..

Wallace, T. and Hussain, A. (1969) "The Use of Error Components Models in Combining Cross Section with Time Series Data", Econometrica, 37, 55-72.

Zellner, A. (1962) "An Efficient Method of Estimating Seemingly Unrelated Regressions and Tests for Aggregation Bias, Journal of the American Statistical Association", 57, 348-368.

A. Appendix

This appendix summarizes the technical foundations and implications of the results presented. It begins with consideration of the two-wave model.

A.1 Two-Wave Model

The first result (Theorem 1) follows directly from the theory of multivariate linear models because the two-wave model with fixed initial wave can be treated as a multivariate regression model. This theory also implies that the distribution of $\{\sqrt{n}(\hat{b}_j - b_j)\}$ is normal for samples of all sizes and thus the normal is not used as a large-sample approximation.

The second result (Theorem 2) also follows from the theory of multivariate linear models since it is a combined application of the multivariate extension of the Gauss-Markov Theorem and the multivariate central limit theorem (Anderson (1984), Rao (1973)).

The third result (Theorem 3) does not follow from the theory of multivariate linear models since the proposition stated is not conditional on the initial wave. The part of the result for the model with normal errors does follow, however, from the theory of replicated vector-valued autoregressive processes developed by Anderson (1978) or by direct calculation. The log of the likelihood function for the model with normal initial observations and normal errors is

$$\begin{aligned}\ell = c & - \frac{1}{2} \log |\theta_0| - \frac{1}{2} \log |\beta| - \frac{n}{2} \text{tr}[\theta_0^{-1} z_0' z_0] \\ & + \text{tr} \beta^{-1} [z_1' z_1 - 2z_0' z_0 \beta' + \beta z_0 z_0 \beta']\end{aligned}$$

which is minimized by letting $\hat{\beta} = \hat{\beta}$ as defined in (3.1).

The remainder of Theorem 3, the large sample distribution for the vector \hat{b}_j of estimated regression parameters for the model without normal errors, follows from applying the multivariate central limit theorem and direct calculation of the covariance matrices.

The assumption of a stationary covariance matrix is more attractive in time-series analysis than in panel analysis. For a simple time series of length $T + 1$ the assumption that θ_0 is of a particular form is an assumption on a single observation, Z_0 . In a $T + 1$ wave panel model the assumption is on all the n observations at wave 0. Should the data not be consistent with the assumption, the covariance matrix θ_0 may vary significantly from θ_t ($t = 1, \dots, T$). For this case, the assumption of stationarity imposes significantly on the estimates of the matrices B and θ_t .

In our notation the standard hypothesis of multivariate regression model is

$$H_0(A, B): \quad ABC = 0$$

where B is the matrix of regression parameters and A and C are contrast matrices, matrices which "pick out" rows and columns (Anderson (1984), Duntzman (1984)). There are no matrices A, B which will give a null hypothesis $H_0(A, B)$ identical to $H_0^{(2)}$: $B_{12} = B'_{21} = 0$.

The theory of likelihood ratio methods is a central theme in statistical theory. This theory yields the large sample distribution of the test statistic and certain optimality properties for the test. (eg. Cox and Hinkley (1974)). Most critical for the model at hand the likelihood ratio test is asymptotically locally most powerful under relatively weak conditions, conditions that obtain for the cross-lagged panel model (eg. Cox and Hinkley (1974)).

The likelihood ratio test statistics for the significance of the cross-effects in the two-wave model (t_1) can be approximated by application of Zellner's seemingly unrelated regression. For the simple case ($p = q = 1$) note that when $H_0^{(2)}$ is true the model can be expressed

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} x_0 & 0 \\ 0 & y_0 \end{bmatrix} \begin{bmatrix} \alpha \\ \epsilon \end{bmatrix} + \begin{bmatrix} f_1 \\ g_1 \end{bmatrix} \quad (A.1)$$

which can be relabeled as

$$u_1 = U_0 \theta + v_1 \quad (A.2)$$

where u_1 and v_1 are $2n$ -vectors, U_0 is a $2n \times 2$ matrix, $\theta = (\alpha, \varepsilon)'$ is a 2-vector, and v_1 has covariance matrix $I_n \otimes \Sigma$ where Σ is a 2×2 matrix.

The maximum likelihood estimators for the model in (A.2) can be approximated as follows. [Zellner (1962); Schmidt (1976); Malinvaud (1980)].

If Σ were known then the maximum likelihood estimator of θ would be the weighted least-squares estimator

$$\hat{\theta} = [U_0' (I_n \otimes \Sigma)^{-1} U_0]^{-1} (I_n \otimes \Sigma)^{-1} u_1 = (\hat{\alpha}, \hat{\varepsilon})' \quad (\text{A.3})$$

and thus the maximum likelihood estimator of B would be

$$\hat{B} = \text{diag}(\hat{\alpha}, \hat{\varepsilon}) = \begin{bmatrix} \hat{\alpha} & 0 \\ 0 & \hat{\varepsilon} \end{bmatrix}$$

Since Σ is unknown \hat{B} is approximated by substituting a consistent estimator for Σ in (A.3). To this end let

$$\hat{\theta} = (U_1' U_0)^{-1} U_0' u_1 = (\hat{\alpha}, \hat{\varepsilon})'$$

be the least squares estimator of θ and let

$$\hat{B}_r = \text{diag}(\hat{\alpha}, \hat{\varepsilon})$$

be the restricted least squares estimator of B . Estimate Σ by

$$\hat{\Sigma}_r = (z_1 - z_0 \hat{B}_r)' (z_1 - z_0 \hat{B}_r) / (n-2)$$

and then approximate the maximum likelihood estimator $\hat{\theta}$ by $\hat{\theta}^* = (\alpha^*, \varepsilon^*)'$ where

$$\hat{\theta}^* = [U_0' (I_n \otimes \hat{\Sigma}_r^{-1}) U_0]^{-1} (I_n \otimes \hat{\Sigma}_r^{-1}) u_1$$

and, in turn, approximate the maximum likelihood estimator $\hat{\beta}_v$ by

$$\hat{\beta}_v^* = \text{diag}(\alpha^*, \beta^*) = \begin{bmatrix} \alpha^* & 0 \\ 0 & \beta^* \end{bmatrix} \quad (\text{A.4})$$

For an estimator of the covariance matrix $\hat{\Sigma}_v$ we use

$$\hat{\Sigma}_v^* = (\hat{\zeta}_1 - \hat{\zeta}_0 \hat{\beta}_v^*)' (\hat{\zeta}_1 - \hat{\zeta}_0 \hat{\beta}_v^*) / (n-2)$$

The theory of seemingly unrelated regression [Zellner (1962), Schmidt (1976)] yields that for a two-wave panel study with fixed initial observations and normal errors if the hypothesis of no cross-effects holds the estimator $\hat{\beta}_v^*$ is asymptotically equivalent to the maximum likelihood estimator $\hat{\beta}_v$ and the associated estimator $\hat{\Sigma}_v^*$ is a consistent estimator of the covariance matrix $\hat{\Sigma}_v$.

The likelihood ratio criterion for testing the hypothesis of no cross-effect is approximated by $t_1^* = -2 \log L_1^*$ where

$$L_1^* = \frac{\left| \hat{\Sigma}_v \right|^{\frac{n}{2}} \exp \left\{ -\frac{1}{2} \text{tr}(\hat{\zeta}_1 - \hat{\zeta}_0 \hat{\beta}_v^*) (\hat{\Sigma}_v^*)^{-1} (\hat{\zeta}_1 - \hat{\zeta}_0 \hat{\beta}_v^*)' \right\}}{\left| \hat{\Sigma}_v^* \right|^{\frac{n}{2}} \exp \left\{ -\frac{1}{2} \text{tr}(\hat{\zeta}_1 - \hat{\zeta}_0 \hat{\beta}_v) \hat{\Sigma}_v^{-1} (\hat{\zeta}_1 - \hat{\zeta}_0 \hat{\beta}_v)' \right\}}$$

For a two-wave panel study with fixed initial observations and normal errors if $p = q = 1$ and the hypothesis of no-cross effect holds then the statistic t_1^* is asymptotically chi-square with two degrees of freedom and the test based on t_1^* is asymptotically equivalent to the likelihood ratio test (based on t_1) and is thus an asymptotically locally most powerful test.

The proof follows from the consistency of the estimators $\hat{\Sigma}^*$ and $\hat{\Sigma}$ and from the theory of likelihood ratio tests and asymptotically equivalent tests (Cox and Hinkley (1974)).

For the more general two-wave panel study (p and q not restricted) the model is replaced by a model of seemingly unrelated multivariate regressions. The construction above extends to this more general case with the test statistic $t_1 = -2 \log L_1$ being asymptotically chi-square with $2pq$ degrees of freedom if $H_0^{(2)}$ is true.

A.2 The Multi-wave Model

The first result for the multi-wave panel model, (Theorem 4) is proved by combining the major result on estimating the parameters in a multivariate autoregressive process [Anderson (1978)] or a result in econometrics on estimating the parameters in a recursive linear system (Malinvaud (1980)) with the multivariate central limit theory (Anderson 1984)).

For the multi-wave model the hypothesis of no cross-effects is the same as for the two-wave model $H_0^{(2)}$: $\beta_{12} = \beta_{21} = 0$. Under this hypothesis the multi-wave panel model can be expressed as

$$\zeta_t = z_{t-1} \begin{bmatrix} \beta_{11} & 0 \\ 0 & \beta_{22} \end{bmatrix} + \varepsilon_t \quad (t = 1, \dots, T)$$

Since the regressor variables are either fixed or lagged endogenous, the maximum likelihood estimator of β can be approximated by extending the scheme based on Zellner's seemingly unrelated regression model. This scheme yields an estimator

$$\hat{\beta}^* = \begin{bmatrix} \hat{\beta}_{11}^* & 0 \\ 0 & \hat{\beta}_{22}^* \end{bmatrix}$$

which is an extension the estimator displayed in (A.4) and which is asymptotically equivalent to the restricted maximum likelihood estimator of β .

Using notation familiar in the analysis of multivariate linear models (cf., Anderson (1984), Rao (1973)) the error "sum of squares" matrix is defined by

$$\xi = n(T-1) \sum_{t=1}^T (\hat{Z}_t - \hat{Z}_{t-1}\hat{B})(\hat{Z}_t - \hat{Z}_{t-1}\hat{B})'$$

and the hypothesis "sum of squares" matrix is defined by

$$\eta = nT(\hat{B}^* - \hat{B})\hat{Z}_{t-1}'\hat{Z}_{t-1}(\hat{B}^* - \hat{B})'$$

The test statistic T_3^* is defined by

$$t_3^* = |\xi| / |\xi + \eta|$$

The theory of likelihood ratio statistics yields that for a multi-wave panel study with fixed initial observations and normal errors if the hypothesis of no cross-effects holds then the test statistic t_3^* is asymptotically chi-square with $2pq$ degrees of freedom. Furthermore the test based on t_3^* is asymptotically equivalent to the likelihood ratio test, based on t_3^* .

The final test of the hypothesis of no cross-effects is a Wald-type test and follows from the asymptotic distribution of the maximum likelihood estimator given in Theorem 4.

Let π be the $2pq$ vector of all elements in B_{12} and B_{21} and let $\tilde{\pi}$ be a $2pq$ vector of the elements in \hat{B}_{12} and \hat{B}_{21} arranged conformally with π . Let $\hat{\Theta}_{\pi\pi}$ be the submatrix of $\hat{\Theta}$ that is the covariance matrix of $\tilde{\pi}$.

Let $\hat{\theta}_{\pi}^*$ be the corresponding submatrix of $\hat{\theta}$; then let

$$t_4 = \begin{pmatrix} \hat{\pi}' \\ \hat{\theta}_{\pi}^{*-1} \\ \hat{\pi} \end{pmatrix}$$

For a multi-wave panel study if the initial observations are fixed and the errors are normal and if the hypothesis of no cross-effects holds, t_4 is asymptotically chi-square with $2pq$ degrees of freedom and the test based on t_4 is asymptotically equivalent to the likelihood ratio test.

TECHNICAL REPORTS
U.S. ARMY RESEARCH OFFICE - CONTRACT DAAG29-82-K-0156

1. "Maximum Likelihood Estimators and Likelihood Ratio Criteria for Multivariate Elliptically Contoured Distributions," T. W. Anderson and Kai-Tai Fang, September 1982.
2. "A Review and Some Extensions of Takemura's Generalizations of Cochran's Theorem," George P.H. Styan, September 1982.
3. "Some Further Applications of Finite Difference Operators," Kai-Tai Fang, September 1982.
4. "Rank Additivity and Matrix Polynomials," George P.H. Styan and Akimichi Takemura, September 1982.
5. "The Problem of Selecting a Given Number of Representative Points in a Normal Population and a Generalized Mills' Ratio," Kai-Tai Fang and Shu-Dong He, October 1982.
6. "Tensor Analysis of ANOVA Decomposition," Akimichi Takemura, November 1982.
7. "A Statistical Approach to Zonal Polynomials," Akimichi Takemura, January 1983.
8. "Orthogonal Expansion of Quantile Function and Components of the Shapiro-Francia Statistic," Akimichi Takemura, April 1983.
9. "An Orthogonally Invariant Minimax Estimator of the Covariance Matrix of a Multivariate Normal Population," Akimichi Takemura, April 1983.
10. "Relationships Among Classes of Spherical Matrix Distributions," Kai-Tai Fang and Han-Feng Chen, April 1984.
11. "A Generalization of Autocorrelation and Partial Autocorrelation Functions Useful for Identification of ARMA(p,q) Processes," Akimichi Takemura, May 1984.
12. "Methods and Applications of Time Series Analysis Part II: Linear Stochastic Models," T. W. Anderson and N. D. Singpurwalla, October 1984.
13. "Why Do Noninvertible Estimated Moving Averages Occur?" T. W. Anderson and Akimichi Takemura, November 1984.
14. "Invariant Tests and Likelihood Ratio Tests for Multivariate Elliptically Contoured Distributions," Huang Hsu, May 1985.
15. "Statistical Inferences in Cross-lagged Panel Studies," Lawrence S. Mayer, November 1985.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|--|------------------------------|--|
| 1. REPORT NUMBER 15 | 2. GOVT ACCESSION NO. N/A | 3. RECIPIENT'S CATALOG NUMBER N/A |
| 4. TITLE (and Subtitle) Statistical Inferences in Cross-Lagged Panel Studies | | 5. TYPE OF REPORT & PERIOD COVERED Technical Report |
| | | 6. PERFORMING ORG. REPORT NUMBER --- |
| 7. AUTHOR(s) Lawrence S. Mayer | | 8. CONTRACT OR GRANT NUMBER(s) DAAG29-85-K-0239 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Statistics, Sequoia Hall Stanford University Stanford, CA 94305 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS P-19065-M |
| 11. CONTROLLING OFFICE NAME AND ADDRESS U. S. Army Research Office Post Office Box 12211 Research Triangle Park NC 27709 | | 12. REPORT DATE November 1985 |
| 14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) | | 13. NUMBER OF PAGES 45pp. |
| | | 15. SECURITY CLASS. (of this report) Unclassified |
| | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |
| 16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited. | | |
| 17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report) NA | | |
| 18. SUPPLEMENTARY NOTES The view, opinions, and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy, or decision, unless so designated by other documentation. | | |
| 19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Panel analysis; regression; autoregressive processes; multivariate regression. | | |
| 20. ABSTRACT (Continue on reverse side if necessary and identify by block number) (See reverse side.) | | |

20. Abstract

Cross-lagged panel studies are statistical studies in which two or more variables are measured for a large number of subjects at each of several waves or points in time. The variable divide naturally into two sets and the primary purpose of the analysis is to estimate and test the cross-effects between the two sets. Such studies are found in the mainstreams of social, behavioral and business research. One approach to this analysis is to express the cross-effects as parameters in regression equations and then use regression methods to estimate and test the parameters. We contribute to this approach by extending the regression model to a multivariate model that captures the correlation between the dependent variables. We develop estimators for the parameters of this model and hypothesis tests for assessing the presence of effects and cross-effects. We demonstrate our results with the analysis of a cross-lagged panel study of the perceptions and attitudes of patients toward a health maintenance organization.